

A Distance Measure for the Analysis of Polar Opinion Dynamics in Social Networks

Victor Amelkin
University of California
Santa Barbara, CA
victor@cs.ucsb.edu

Ambuj K. Singh
University of California
Santa Barbara, CA
ambuj@cs.ucsb.edu

Petko Bogdanov
University at Albany – SUNY
Albany, NY
pbogdanov@albany.edu

ABSTRACT

Analysis of opinion dynamics in social networks plays an important role in today's life. For applications such as predicting users' political preference, it is particularly important to be able to analyze the dynamics of competing opinions. *While observing the evolution of polar opinions of a social network's users over time, can we tell when the network "behaved" abnormally? Furthermore, can we predict how the opinions of the users will change in the future? Do opinions evolve according to existing network opinion dynamics models?* To answer such questions, it is not sufficient to study individual user behavior, since opinions can spread far beyond users' egonets. We need a method to analyze opinion dynamics of all network users simultaneously and capture the effect of individuals' behavior on the global evolution pattern of the social network.

In this work, we introduce Social Network Distance (SND)—a distance measure that quantifies the "cost" of evolution of one snapshot of a social network into another snapshot under various models of polar opinion propagation. SND has a rich semantics of a transportation problem, yet, is computable in time linear in the number of users, which makes SND applicable to the analysis of large-scale online social networks. In our experiments with synthetic and real-world Twitter data, we demonstrate the utility of our distance measure for anomalous event detection. It achieves a true positive rate of 0.83, twice as high as that of alternatives. When employed for opinion prediction in Twitter, our method's accuracy is 75.63%, which is 7.5% higher than that of the next best method.

Code: <http://cs.ucsb.edu/~victor/pub/ucsb/dbl/snd/>

1. INTRODUCTION

Analysis of people's opinions plays an important role in today's life, and social networks provide a great platform for such research. Businesses are interested in advertising their products in social networks relying on viral marketing. Political strategists are interested in predicting an election outcome based on the observed sentiment change of a sample of voters. Mass media and security analysts may be interested in a timely discovery of anomalies based on how a social network "behaves". Thus, it is important to have methods for the analysis of how user opinions evolve in a social network.

How can we quantify the change in opinions of users with respect to their expected behavior in a social network? Can we distinguish opinion shifts caused by in-network user in-

teraction from those caused by factors external to the network? To answer such questions, we need a distance measure that explicitly models opinion dynamics, incorporating both the distribution of user opinions at two time instances and the network structure that defines the pathways for opinion dissemination. In this work, we develop such a distance measure for snapshots of a social network and employ it for the analysis of competing opinion dynamics.

While the dynamics of a social network can be characterized by evolution of both the network's structure and user opinions [23], in this paper we focus on the opinion dynamics. We assume that there are two *polar opinions* in the network, *positive* "+" and *negative* "-". Users having no or an unknown opinion are termed *neutral*, while those expressing opinion—*active*. A *network state* is comprised of the opinions of all network users at a given time. Polar opinions *compete* in that a user having a positive opinion is unlikely to enthusiastically spread information about the adverse negative opinion, yet, would spread information about the friendly positive opinion "at a discount cost". Such competition may arise when the notions the opinions relate to are inherently competing. For example, in an election, the voters leaning toward one political party are unlikely to spread positive rumors about the competing party. Another example is viral marketing, where the consumers who favor smartphones of one brand may readily express their affection to it, but may be unwilling to praise the competing brand.

Given a time series of network states, we address the applications of detecting anomalous network states and predicting opinions of individual users. For the first application, we answer the question of which network states are anomalous with respect to the observed evolution of the network. For the second application, we predict currently unknown opinions of selected users in the network based on the network's observed past and present behavior.

The analysis of a time series of network states is, however, complicated, because network states do not naturally belong to any vector space, and the existing time series analysis techniques cannot be readily applied. Our approach is to treat network states as members of a metric space induced by a distance measure governed by both the network's structure and user opinions. We design a semantically and mathematically appealing as well as efficiently computable distance measure *Social Network Distance (SND)* for the comparison of social network states containing polar opinions, and demonstrate its applicability to the analysis of real-world data.

To quantify the dissimilarity between network states, SND

takes into account how information propagates in the network. A change of a given user’s opinion from, say, neutral to positive, contributes to the overall distance between the corresponding network states by reflecting the likelihood of this user’s opinion change based on the opinions and locations of other users in the network under a chosen model of polar opinion dynamics. However, since the network users interact, the distance measure needs to consider the opinion shifts of all users simultaneously. Thus, we define SND as a transportation problem that models opinion spread and adoption in the network. In particular, by making the transportation costs dependent on both the network’s topology and the opinions of the users conducting information in the network, we capture the competitive aspect of polar opinion propagation.

The summary of *our contributions* is as follows:

- We propose SND—the first distance measure suitable for comparison of social network states containing competing opinions under various models of opinion dynamics.
- We develop a scalable method to precisely compute SND in time linear in the number of users in the network, thus, making SND applicable to the analysis of real-world online social networks.
- We demonstrate the applicability of our distance measure using both synthetic and real-world data from Twitter. In detecting anomalous states of a social network, SND is superior to other distance measures in discovering controversial events that have polarized the society. In user opinion prediction experiments, SND also outperforms competitors in terms of prediction accuracy.

2. EARTH MOVER’S DISTANCE AND NETWORK STATE COMPARISON

A good distance measure for the states of a social network should take into account the specifics of polar opinion propagation in a network. For example, a user having opinion “+” should not actively participate in the dissemination of opinion “−”, or, at least, the dissemination of an adverse opinion should incur a large cost. On the other hand, this user should disseminate friendly opinion “+” at a cost lower than the cost a neutral user would incur. Thus, we propose to address the problem of comparing states of a social network as a transportation problem where the costs of opinion propagation are defined based on the shortest paths between the users in the network, computed taking into account both the network structure and the user opinions.

This high-level intuition about network state comparison as an opinion transportation problem inadvertently leads us to one of the well-studied distance measures—Earth Mover’s Distance (EMD). Originally, defined as a dissimilarity measure for histograms [25], EMD can be used for the comparison of network states viewed as histograms, with histogram bins’ values quantifying individual user opinions. Intuitively, EMD measures the costs of optimal transformation of one histogram into another with respect to the *ground distance* specifying the costs of moving mass between bins. In our case, the ground distance is defined based on the shortest paths between the users of the network, where the shortest paths depend on both the network’s topology as well as the opinions of the users facilitating opinion propagation.

Formally, given two real-valued histograms $P = [P_1, \dots, P_n]$ and $Q = [Q_1, \dots, Q_m]$, EMD between them with respect to a cross-bin ground distance $\{D_{ij}\}_{n \times m}$ is defined as the so-

lution to the problem of optimal mass transportation from the bins of P (suppliers) to the bins of Q (consumers) with respect to transportation costs D .

$$\text{EMD}(P, Q, D) = \sum_{i=1}^n \sum_{j=1}^m D_{ij} \hat{f}_{ij} / \sum_{i=1}^n \sum_{j=1}^m \hat{f}_{ij}, \quad (1)$$

where $\{\hat{f}_{ij}\}_{n \times m}$ is an optimal transportation plan in the following transportation problem:

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^m f_{ij} D_{ij} \rightarrow \min, \quad \sum_{i=1}^n \sum_{j=1}^m f_{ij} &= \min \left\{ \sum_{i=1}^n P_i, \sum_{j=1}^m Q_j \right\}, \\ f_{ij} \geq 0, \quad \sum_{j=1}^m f_{ij} &\leq P_i, \quad \sum_{i=1}^n f_{ij} \leq Q_j, \quad (1 \leq i \leq n, 1 \leq j \leq m). \end{aligned}$$

EMD is attractive not only because its ground distance can capture how opinions propagate in the underlying network, but it is also driven by node states rather than network topology, is spatially-sensitive, and metric under the following conditions.

THEOREM 1. (Rubner et al. [25]) *If all histograms under comparison have equal total masses, and the underlying ground distance is metric, then EMD is metric.*

In the following section, we use EMD to construct a distance measure for network states containing polar opinions.

3. DISTANCE MEASURE FOR NETWORK STATES WITH POLAR OPINIONS

Given a network $G = \langle V, E \rangle$, where V ($|V| = n$) is the set of nodes (users) and E is the set of edges (social ties), we want to compute the distance between two of its states $P = [P_1, \dots, P_n]^T$ and $Q = [Q_1, \dots, Q_n]^T$. While generalizations are possible, we use an intuitive and simple scheme for polar opinion quantification: if user i has opinion “+” in network state P , then $P_i = +1$; $P_i = -1$ if the user’s opinion is “−”; and $P_i = 0$ if the user is neutral¹.

Despite the appeal of EMD, it is not readily applicable to the comparison of network states P and Q , since (i) the users’ behavior may change in the process of opinion propagation, while a transportation problem underlying EMD operates with static transportation costs; (ii) EMD is defined for histograms of a homogeneous quantity, while P and Q contain both positive and negative values; and (iii) it is not clear how to incorporate the node states into the definition of the ground distance. In order to define our Social Network Distance (SND), we address these three problems in what follows.

(i) SND as a transportation problem. For two given users u and v , not necessarily being immediate neighbors in the network, the cost of v ’s adopting opinion from u depends not only on the states and locations of u and v in the network, but also on the states of the users through which u ’s opinion can reach v .

In Fig. 1.a, user v_3 having opinion “+” affects the cost of propagating opinion “−” from user v_1 to user v_4 . In Fig. 1.b, however, user v_3 is initially neutral, but can become active in

¹ There is a great body of research on methods for opinion classification based on user-generated content. In this work, however, we assume that we have access only to the quantified opinions, and no access to the user-generated content (e.g., tweets), which may be unavailable due to privacy reasons.

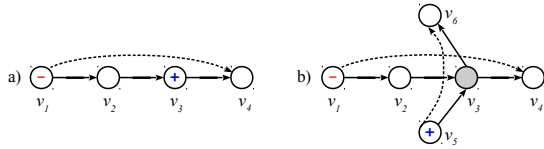


Figure 1: Transitive opinion propagation in a social network. The solid and dashed arrows represent social ties and opinion flow, respectively.

the process of v_5 's propagating "+" to v_6 before opinion "-" from v_1 reaches v_3 , thereby, impeding the spread of "-" from v_1 to v_4 . In order to pose SND as a transportation problem, we will assume that the costs of opinion propagation depend only on the opinions of the currently active users, taking no account of the potential change of user opinions in the process of opinion transportation.

(ii) Handling polar opinions. We design SND to measure the optimal cost of transforming one network state G_1 into another network state G_2 by the means of opinion transportation. The active users of G_1 are the suppliers and those in G_2 are the consumers in the transportation problem setting. In defining the transportation costs, we assume that users adopting, say, opinion "+", in G_2 are affected by others of the same opinion during the propagation. Similarly, suppliers only propagate opinions to the consumers of the same type. Thus, the set of constraints in the transportation problem can be divided into two non-overlapping subsets for two kinds of opinions. Consequently, the problem of optimal transportation of opinions from suppliers of G_1 to consumers of G_2 can be split into two transportation problems: one for transporting the opinions of each kind.

(iii) Defining the ground distance. The cost of opinion propagation from user u to user v depends on their topological proximity, how frequently they communicate, persuasiveness, and stubbornness of u and v as well as the users "separating" them. Formally, the ground distance $D(G_i, op) \in \mathbb{R}^{+n \times n}$, reflecting the costs of propagating opinion op through a network in state G_i , is a matrix containing the lengths of the shortest paths in a network with its adjacency matrix defined as:

$$A^{ext}(G_i, op) = -\log \mathbb{P}(G_i, op) - \log \mathbb{P}^{in}(G_i, op) - \log \mathbb{P}^{out}(G_i, op), \quad (2)$$

where the summands on the right are n -by- n matrices of log-probabilities of communication, opinion adoption, and opinion spreading, respectively. Probabilities $\mathbb{P}(G_i, op)$ can be defined as the relative frequencies of communication between users. In the absence of such information, we set $-\log \mathbb{P}(G_i, op)$ to be the connectivity matrix of the network, penalizing for the users' topological remoteness. Opinion adoption probabilities $\mathbb{P}^{in}(G_i, op)$ reflect users' susceptibility/stubbornness [28]. If such information is unavailable, for each existing edge $\langle u, v \rangle$, we set $\mathbb{P}_{uv}^{in} = 1$, so that all users are non-stubborn and equally receptive to persuasion. Finally, the opinion spreading penalties $-\log \mathbb{P}^{out}(G_i, op)$ are defined based on a particular opinion dynamics model. Several ways to make such an assignment are described below.

Model-agnostic Opinion Propagation: If there is no evidence that opinions evolve in the network according to a specific opinion dynamics model, then the opinion spread-

ing penalties can be defined as

$$-\log \mathbb{P}_{uv}^{out}(G_i, op) = \begin{cases} c_{adverse} & \text{if } G_i[u] \neq op \vee G_i[v] = -op, \\ c_{neutral} & \text{if } G_i[u] = 0, \\ c_{friendly} & \text{if } G_i[u] = op, \end{cases}$$

where $c_{adverse}$, $c_{neutral}$, $c_{friendly}$ are constant penalties for spreading opinion op by the users having respectively adverse, neutral, or friendly opinion relative to op , and $G_i[u]$ is the opinion of user u in network state G_i . This simple definition implies that users willingly spread opinions similar to their own ($c_{friendly}$ is small); are unwilling to spread adverse opinions ($c_{adverse}$ is large); with neutral users' behavior being somewhere in-between ($c_{friendly} < c_{neutral} < c_{adverse}$).

Alternatively, $\mathbb{P}_{uv}^{out}(G_i, op)$ can be defined via one of the existing opinion dynamics models discussed next.

Independent Cascade Model: The distance-based model of Carnes et al. [7] is a version of Independent Cascade Model capturing opinion competition. In this model, we have two sets of initial adopters I_+ , I_- of opinions "+" and "-", respectively, with $I = I_+ \cup I_-$. Each edge $\langle u, v \rangle$ is labeled with an activation probability p_{uv} (which can be learned from the observed data [13]) and a distance d_{uv} . If we denote by $d_v(I)$ the shortest distance from any user of set I to user v with respect to edge distances d_{uv} , and denote by $p^a(G_i, v)$ the sum of edge activation probabilities p_{uv} taken over all active users u in G_i such that $d_v(u) = d_v(I)$, then

$$\mathbb{P}_{uv}^{out}(G_i, op) = \begin{cases} 0 & \text{if } d_v(\{u\}) > d_v(I), \\ 1 & \text{if } G_i[u] = op \wedge G_i[v] = op, \\ \frac{\max(0, p_{uv} - \epsilon)}{p^a(G_i, v)} & \text{if } G_i[u] = op \wedge G_i[v] = 0, \\ \epsilon & \text{otherwise.} \end{cases}$$

In the original model of [7], $\epsilon = 0$, that is, neutral users cannot infect others, and active users do not drop their opinions or spread competing opinions. If, however, we compare network states with respect to the original model, then many network states derived from real-world data would be at distance $+\infty$ from each other, either due to the lack of knowledge about the network (an edge has not been observed or a user's opinion has been misclassified) or due to an imperfect fit of the model and the data. Instead of just declaring two network states as qualitatively unreachable, we always want to quantify the distance between them, and, thus, assign some negligible probabilities ϵ to the events that the model posits as impossible.

Linear Threshold Model: A version of Linear Threshold Model supporting opinion competition has been proposed by Borodin et al. [5]. In this model, each edge $\langle u, v \rangle$ is weighted with ω_{uv} reflecting the amount of influence u has over v ; and each user u has an in-advance chosen constant threshold θ_u . If we denote by $N^{in}(G_i, v)$ the set of in-neighbors of v active in G_i , and by Ω^{in} the sum of ω_{xv} over all $x \in N^{in}(G_i, v)$, then

$$\mathbb{P}_{uv}^{out}(G_i, op) = \begin{cases} 0 & \text{if } u \notin N^{in}(G_i, v), \\ 1 & \text{if } G_i[u] = op \wedge G_i[v] = op, \\ \frac{(1-\epsilon)\omega_{uv}}{\Omega^{in}} & \text{if } G_i[u] = op \wedge G_i[v] = 0 \wedge \Omega^{in} \geq \theta_v, \\ \epsilon & \text{otherwise.} \end{cases}$$

Having addressed challenges (i)-(iii), we are now ready to formally define SND.

In general, the network's structure might have changed between the times corresponding to the two network states under comparison [23], but defining the ground distance for each pair of users over a different network would incur an

unacceptably high time complexity of the resulting distance measure. Thus, for time-ordered network states G_1 and G_2 , one can define the ground distance based on the network structure corresponding to the earlier network state G_1 . However, to make SND applicable to time-unordered network states as well, we define SND based on both $D(G_1, op)$ and $D(G_2, op)$ as follows.

$$\begin{aligned} \text{SND}(G_1, G_2) &= \frac{1}{2} \times [\\ &\text{EMD}(G_1^+, G_2^+, D(G_1, +)) + \text{EMD}(G_1^-, G_2^-, D(G_1, -)) + \\ &\text{EMD}(G_2^+, G_1^+, D(G_2, +)) + \text{EMD}(G_2^-, G_1^-, D(G_2, -))], \end{aligned} \quad (3)$$

where users having opinion “−” are considered neutral in G_i^+ , users having opinion “+” are neutral in G_i^- , and EMD is the original Earth Mover’s Distance [25] that will further be replaced with its generalization $\widehat{\text{EMD}}$ designed in the next section. Since SND is a linear combination of several instances of EMD, then SND preserves most mathematical properties, and, in particular, metricity of the chosen EMD.

4. GENERALIZED EARTH MOVER’S DISTANCE

The original EMD [25] is limited in that it cannot adequately compare histograms having different total mass—it ignores the mass mismatch, assigning a small distance value to a pair of a light and a heavy histograms. However, if we think about two histograms corresponding to the states of a social network, one with a few and another one with many active users, then we expect the distance between such histograms to be large. In real-world data, even consecutive network states observe a widely varying number of active users, making the challenge of comparing histograms with total mass mismatch well pronounced.

There are several extensions of EMD that address the above mentioned limitation. One of them, $\widehat{\text{EMD}}$ [24], augments EMD with an additive mass mismatch penalty as follows

$$\begin{aligned} \widehat{\text{EMD}}(P, Q, D) &= \text{EMD}(P, Q, D) \cdot \min \left\{ \sum P_i, \sum Q_j \right\} + \\ &+ \alpha \cdot \max \{D_{ij}\} \cdot \left| \sum P_i - \sum Q_j \right|, \end{aligned}$$

where EMD is the original Earth Mover’s Distance, and α is a constant parameter. The second summand represents the mass mismatch penalty that depends only on the magnitude of the mass mismatch and the maximum ground distance, thereby, being unable to capture the fine details of the network’s structure D can depend upon. This is, however, inadequate for the comparison of the states of a social network, because the network’s behavior depends not only on the number of new activations, but also on where these newly activated users are located in the network.

Another EMD version, namely, EMD^α [18], extends each histogram with an extra bin (“the bank bin”) whose value is chosen so that the total masses of the histograms become equal. An example of such an extension is shown in Fig. 2. Formally, EMD^α is defined as follows.

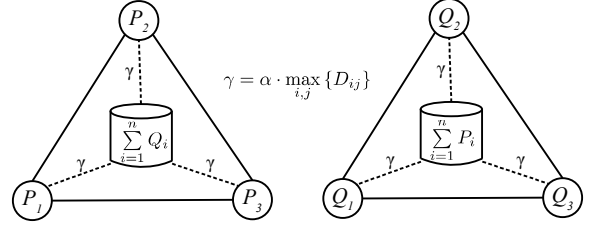


Figure 2: Histograms P and Q defined over the same network are extended with bank bins, whose capacities are chosen, so that the total masses of the extended histograms \tilde{P} and \tilde{Q} are equal. The ground distances $\tilde{D}_{bank,i} = \tilde{D}_{i,bank} = \gamma$ from and to the bank bin are uniformly defined based on the largest ground distance between the initially present bins.

$$P = [P_1, \dots, P_n], \quad Q = [Q_1, \dots, Q_n],$$

$$P_{bank} = \sum_{j=1}^n Q_j, \quad \tilde{P} = [P, P_{n+1} = P_{bank}],$$

$$Q_{bank} = \sum_{i=1}^n P_i, \quad \tilde{Q} = [Q, Q_{n+1} = Q_{bank}],$$

$$\tilde{D} = \left[\begin{array}{c|c} D_{n \times n} & \alpha \max_{i,j} \{D_{ij}\} \\ \hline -\alpha \max_{i,j} \{D_{ij}\} - & 0 \end{array} \right],$$

$$\text{EMD}^\alpha(P, Q) = \text{EMD}(\tilde{P}, \tilde{Q}, \tilde{D}) \cdot \left(\sum_{i=1}^n P_i + \sum_{j=1}^n Q_j \right).$$

However, as we establish in Theorem 2, EMD^α is equivalent to $\widehat{\text{EMD}}$ and, hence, is also unsuitable for the comparison of social network states.

THEOREM 2. *If ground distance $D \in \mathbb{R}^{n \times n}$ and parameter $\alpha \in \mathbb{R}^+$ are chosen such that both EMD^α and $\widehat{\text{EMD}}$ are metric, that is, D is metric and $\alpha \geq 0.5$ [18, 24], then $\forall P, Q \in \mathbb{R}^n : \text{EMD}^\alpha(P, Q, D) = \widehat{\text{EMD}}(P, Q, D)$.*

PROOF. W.l.o.g., let us assume that $\sum P_i \leq \sum Q_j$. The proof will use the following notation:

$$\Delta = \Delta(P, Q) = \left| \sum_{i=1}^n P_i - \sum_{j=1}^n Q_j \right|, \quad \gamma = \alpha \max_{i,j} \{D_{ij}\}.$$

Hence, $\widehat{\text{EMD}}$ is defined as

$$\widehat{\text{EMD}}(P, Q) = \text{EMD}(P, Q) \min \left\{ \sum_{i=1}^n P_i, \sum_{j=1}^n Q_j \right\} + \gamma \Delta.$$

The goal of the proof is to show that EMD^α has exactly the same expression as $\widehat{\text{EMD}}$, as long as they are metric. Consider how a unit of mass can be transported between two histograms (see Fig. 3). There are two qualitatively different alternatives for moving a unit of mass from regular (non-bank) bin i of histogram \tilde{P} : a unit of mass can be moved either to a regular bin j or the bank bin of \tilde{Q} .

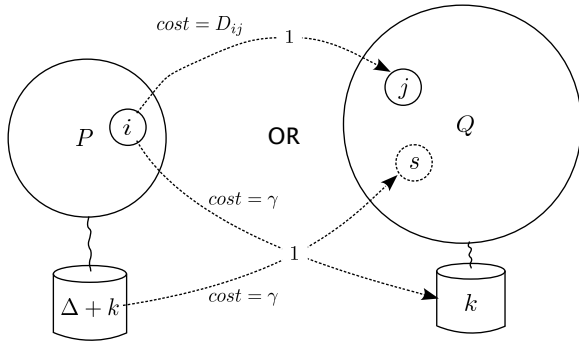


Figure 3: Two qualitatively different ways to transport a unit of mass from extended histogram $\tilde{P} = [P, k + \Delta]$ to extended histogram $\tilde{Q} = [Q, k]$. (Dashed arrows represent the flow of mass.) The bank bin is attached to every node of each histogram. $k = \sum P_i$, so that the total masses of two histograms are equal.

In the first case, the total cost of transportation of a unit of mass is exactly the ground distance $\tilde{D}_{ij} = D_{ij}$ between regular bins i and j .

In the second case, the immediate cost of transportation to the bank bin is $\tilde{D}_{i,bank} = \gamma$. However, because we have routed mass from a regular bin to the bank bin, there exists a regular bin s in \tilde{Q} having “mass deficit” that has to be fulfilled from the bank bin of \tilde{P} . Thus, if we move a unit of mass from a regular bin of \tilde{P} to the bank bin of \tilde{Q} , there is an additional incurred cost γ of moving an additional unit of mass from the bank bin of \tilde{P} to some regular bin of \tilde{Q} . Hence, the total cost of transportation of a unit of mass in the second case is

$$\gamma + \gamma = 2\alpha \max_{i,j} D_{ij} \geq (\text{since } \alpha \geq 0.5) \geq \max_{i,j} D_{ij}.$$

Thus, from the point of view of optimal mass transportation, it may never be preferable to move a unit of mass from a regular bin to the bank bin if there is an option to transport mass from a regular bin to a regular bin. Consequently, an optimal solution to the EMD^α 's transportation problem can be decomposed as follows.

$$\begin{aligned} \text{EMD}^\alpha(P, Q, D) &= \text{EMD}(\tilde{P}, \tilde{Q}, \tilde{D}) \left(\sum_{i=1}^n P_i + \sum_{j=1}^n Q_j \right) \\ &= \min_{\{f_{ij}\}} \sum_{i,j=1}^{n+1} f_{ij} \tilde{D}_{ij} = (\text{let } n+1 = b) \\ &= \min_{\{f_{ij}\}} \left[\underbrace{\sum_{i,j=1}^n f_{ij} \tilde{D}_{ij}}_{\text{regular bins to regular bins}} + \underbrace{\sum_{i=1}^n f_{ib} \tilde{D}_{ib}}_{\text{regular bins to bank bin}} + \underbrace{\sum_{j=1}^n f_{bj} \tilde{D}_{bj}}_{\text{bank bin to regular bins}} + \underbrace{f_{bb} \tilde{D}_{bb}}_{\text{bank bin to bank bin}} \right] \\ &= \min_{\{f_{ij}\}} \left[\sum_{i,j=1}^n f_{ij} D_{ij} + \underbrace{\gamma \sum_{j=1}^n f_{bj}}_{\Delta} \right] = \min_{\{f_{ij}\}} \left[\sum_{i,j=1}^n f_{ij} D_{ij} + \gamma \Delta \right] \\ &= \min_{\{f_{ij}\}} \left[\sum_{i,j=1}^n f_{ij} D_{ij} \right] + \gamma \Delta = \widehat{\text{EMD}}(P, Q, D). \end{aligned}$$

An additional useful observation is that a particular value of

k does not matter for EMD^α , since for every optimal solution of its underlying transportation problem, any amount of mass to the excess of Δ in the bank bin of the lighter histogram will be transported at zero-cost $\tilde{D}_{bank,bank}$ to the bank bin of the heavier histogram. This observation is formalized in the following corollary. \square

COROLLARY 1. *For histograms $P = [P_1, \dots, P_n]$ and $Q = [Q_1, \dots, Q_n]$, and ground distance D , if $\sum P_i = \sum Q_j$ and D is metric, then for all $k \geq 0 \in \mathbb{R}^+$, the following holds.*

$$\text{EMD} \left([P, k], [Q, k], \left[\begin{array}{c|c} D & \omega \\ \hline -\omega & 0 \end{array} \right] \right) = \text{EMD}(P, Q, D),$$

where $[X, k]$ is histogram X extended with a single bank bin with capacity k and a uniformly defined ground distance $\omega \geq \frac{1}{2} \max_{i,j} D_{i,j}$ to/from the regular bins of X . In other words, if two histograms have equal total masses, we can increase their total masses by an arbitrary non-negative k without affecting EMD between the histograms.

Our version of Earth Mover's Distance, EMD^* , extends the idea of EMD^α by augmenting the histograms to even their masses. However, unlike EMD^α , EMD^* extends the histograms with multiple *local bank bins* and distributes the total mass mismatch over all of them. We, thereby, relate the mass mismatch penalty to the network structure, while achieving equality of the total masses of the two histograms under comparison.

With respect to the number and location of the bank bins, one extreme option is to attach one local bank bin to each initially present bin. Furthermore, in order to model a non-constant transportation cost to/from a bank bin, multiple local bank bins per initially present bin can be used, each with its individual ground distance. This bank allocation strategy can incur a high computational cost, since attaching even a single bank bin to each initially present bin doubles the size of the histogram.

A compromise between the two extremes of having a single bank (as in EMD^α) and one bank per bin is to use one or more local banks per a *group of bins*. Such bin groups can be defined based on the structural proximity of the corresponding users in the underlying network (see Fig. 4).

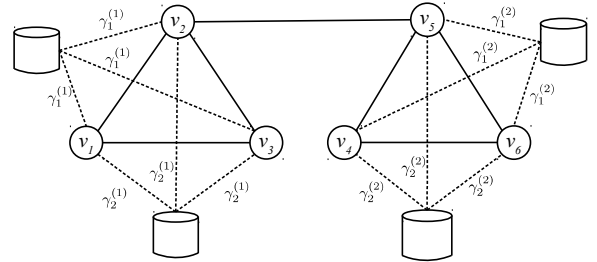


Figure 4: A histogram over a network extended with two banks per cluster of bins. $\gamma_j^{(i)}$ is the ground distance to/from the j 'th bank bin of i 'th bin cluster.

Ground distance $\gamma^{(i)}$ to/from an added bank bin should be of the same order as the ground distances within the i 'th cluster of bins the bank is attached to. If $\gamma^{(i)}$ is much lower

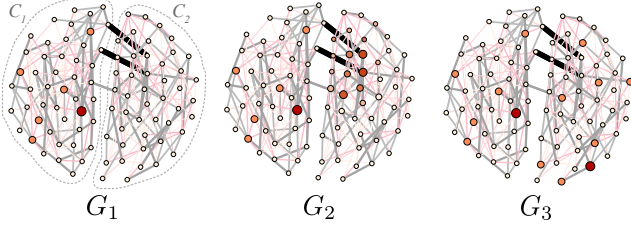


Figure 5: Three histograms over the same network.

than the ground distances in its cluster, then it can negatively affect metricity of EMD^* , the conditions for which will be stated later. If $\gamma^{(i)}$ is much higher than the ground distances in its cluster, it may result in an EMD^α -like behavior, with the global bank bin, even though spatially distributed across multiple local banks, still playing no role in the process of optimal mass transportation. The capacities of the added bank bins should be determined based upon two ideas. Firstly, the capacity of a bank bin should intuitively be proportional to the total mass of the bins the bank is attached to, thereby, preserving the relative distribution of mass over the network. Secondly, the capacities of all the bank bins should be such, that the two histograms under comparison have equal total masses. The following definition of capacity $P^{(i)}$ of a bank bin connected to the i 'th cluster C_i of bins in the context of comparing histograms $P = [P_1, \dots, P_n]$ and $Q = [Q_1, \dots, Q_n]$ incorporates both of the above requirements.

$$P_j^{(i)} = \begin{cases} \sum_{v_k \in C_i} P_k / \left(\sum_{s=1}^n Q_s - \sum_{s=1}^n P_s \right) & \text{if } \sum Q_s > \sum P_s, \\ 0, & \text{otherwise.} \end{cases}$$

To better understand the advantage of EMD^* over the existing versions of EMD, consider the example in Fig. 5. There are three histograms over a network with two pronounced clusters C_1 and C_2 connected by three bridge edges. The distribution of mass over cluster C_1 is identical in all three histograms, while cluster C_2 is empty in G_1 and has some differently distributed mass in G_2 and G_3 . In G_2 the extra mass has been “propagated” from cluster C_1 to cluster C_2 through the bridges, while in G_3 the same amount of extra mass has been randomly distributed over cluster C_2 . Thus, if we assume that G_2 and G_3 have “evolved” from G_1 through a process of mass propagation, then G_2 should intuitively be closer to G_1 than G_3 . However, only EMD^* captures this intuition as $\text{EMD}^*(G_1, G_2) < \text{EMD}^*(G_1, G_3)$, while for EMD^α and $\overline{\text{EMD}}$, G_2 and G_3 are equidistant from G_1 , and for EMD, both G_2 and G_3 are identical to G_1 .

Next, we formally define EMD^* . Suppose we are given two histograms $P = [P_1, \dots, P_n]$ and $Q = [Q_1, \dots, Q_n]$ defined over a network $G = \langle V, E \rangle$ with cross-bin ground distance $D_{n \times n}$. The ground distance is application-specific and, in our case, is provided by SND. Bins of both histograms are clustered into groups $\{C_i\}$, $i = 1, \dots, N_c$ based on the proximity of the corresponding users in network G . Cluster C_i contains NC_i users/bins. Each bin cluster gets N_b banks attached to all its bins, so the total number of bins in an extended histogram is $N = n + N_c \times N_b$. Ground distances from/to the bins of cluster C_i to/from its banks are defined as $\gamma^{(i)} = [\gamma_1^{(i)}, \dots, \gamma_{N_b}^{(i)}]^T$, and, jointly for all clusters, $\gamma = [(\gamma^{(1)})^T, \dots, (\gamma^{(N_c)})^T]^T$. Since bank bins belonging to different clusters are not interconnected, in order to define ground distances between bank bins we define distances between

clusters C_i in terms of D_{ij} as $d_{ij} = \min_{v_p \in C_i, v_q \in C_j} \{D_{pq}\}$.

Then, EMD^* is defined as follows.

$$\tilde{P} = \left[\underbrace{P_1, \dots, P_n}_{\text{original } P}, \underbrace{P_1^{(1)}, \dots, P_{N_b}^{(1)}}_{\text{cluster } C_1 \text{ banks}}, \dots, \underbrace{P_1^{(N_c)}, \dots, P_{N_b}^{(N_c)}}_{\text{cluster } C_{N_c} \text{ banks}} \right],$$

$$\tilde{Q} = \left[\underbrace{Q_1, \dots, Q_n}_{\text{original } Q}, \underbrace{Q_1^{(1)}, \dots, Q_{N_b}^{(1)}}_{\text{cluster } C_1 \text{ banks}}, \dots, \underbrace{Q_1^{(N_c)}, \dots, Q_{N_b}^{(N_c)}}_{\text{cluster } C_{N_c} \text{ banks}} \right],$$

$$S = \left[d_{1,*} \otimes \mathbb{1}_{NC_1 \times 1} \mid \dots \mid d_{N_c,*} \otimes \mathbb{1}_{NC_{N_c} \times 1} \right]^T,$$

$$\tilde{D} = \left[\begin{array}{c|c} D_{n \times n} & \mathbb{1}_{n \times 1} \otimes \gamma^T + S^T \otimes \mathbb{1}_{1 \times N_b} \\ \hline \mathbb{1}_{1 \times n} \otimes \gamma + S \otimes \mathbb{1}_{N_b \times 1} & \gamma \otimes \mathbb{1}_{1 \times (N_b \cdot N_c)} + \gamma^T \otimes \mathbb{1}_{(N_b \cdot N_c) \times 1} - 2 \cdot \text{diag}(\gamma) + d \otimes \mathbb{1}_{N_b \times N_b} \end{array} \right],$$

$$\text{EMD}^*(P, Q) = \text{EMD}(\tilde{P}, \tilde{Q}, \tilde{D}) \max \left\{ \sum P_i, \sum Q_j \right\}, \quad (4)$$

where $\mathbb{1}_{a \times b}$ is an a -by- b matrix of all ones; $d_{*,j}$ and $d_{i,*}$ are the j 'th column and the i 'th row of matrix d , respectively; $\text{diag}(v)$ is a diagonal matrix with the elements of vector v on its main diagonal; and \otimes is the Kronecker product.

Metricity of EMD^* , which can be exploited to improve practical performance of distance-based search in applications [8], is established in the following Theorem 3.

THEOREM 3. *Given an arbitrary finite set \mathcal{H} of histograms with bin clusters $\{C_i\}$ and ground distance $D_{n \times n}$, if D is metric and the ground distances γ to/from the bank bins are such that $\forall i, j : \gamma_j^{(i)} \geq \frac{1}{2} \max_{v_p, v_q \in C_i} D_{pq}$, then EMD^* defined with $\{C_i\}$ and γ is metric on $\mathcal{H} \times \mathcal{H}$.*

PROOF. Let us define constant $M = \max_{X \in \mathcal{H}} \sum_k X_k$. Since \mathcal{H} is finite and all the distributions are assumed to have finite total masses, then $M < +\infty$. Next, we define an auxiliary distance measure EMD' as follows.

$$\text{EMD}'(P, Q, D) = \text{EMD}(P', Q', D'),$$

$$P' = [\tilde{P}, M - \sum \tilde{P}_i], \quad Q' = [\tilde{Q}, M - \sum \tilde{Q}_j],$$

$$D' = \left[\begin{array}{c|c} \tilde{D} & \max_{i,j} \{ \tilde{D}_{ij} \} / 2 \\ \hline - \max_{i,j} \{ \tilde{D}_{ij} \} / 2 & 0 \end{array} \right],$$

where \tilde{P} , \tilde{Q} , and \tilde{D} are the extended histograms, and the extended ground distance, respectively, as defined by EMD^* . From the definition of EMD^* (4), it follows that $\sum \tilde{P}_i = \sum \tilde{Q}_j$ and, hence $M - \sum \tilde{P}_i = M - \sum \tilde{Q}_j = k$. Thus, since $\sum \tilde{P}_i = \sum \tilde{Q}_j = M$, D is metric, and $k \geq 0$, we can apply Corollary 1 to $P = \tilde{P}$, $Q = \tilde{Q}$, $D = \tilde{D}$, and $\omega = \frac{1}{2} \max_{i,j} \tilde{D}_{ij}$, to obtain

$$\begin{aligned} \text{EMD}'(P, Q, D) &= \text{EMD}(P', Q', D') = \\ &= (\text{from Corollary 1}) = \text{EMD}(\tilde{P}, \tilde{Q}, \tilde{D}) = \\ &= (\text{from definition of } \text{EMD}^*) = \frac{\text{EMD}^*(P, Q, D)}{\max \{ \sum P_i, \sum Q_j \}}. \end{aligned}$$

Thus, EMD^* is metric iff EMD' is metric. The latter's metricity, according to Theorem 1, requires equality of total

masses of all histograms and metricity of the ground distance. From the definition of EMD' , it is clear that *all* histograms P' and Q' supplied to EMD by EMD' always have the same total mass M . As to metricity of the ground distance, the identity of indiscernibles and symmetry straightforwardly follow from the corresponding properties of the original ground distance D and our choice of the ground distances to/from the bank bins to be non-negative and symmetric. The triangle inequality holds for the original D , so we need to inspect only the new “triangles” introduced into the ground distance after the addition of the bank bins, such as shown in Fig. 6.

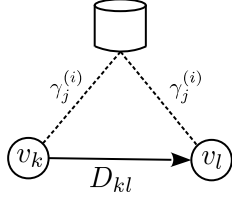


Figure 6: j 'th bank bin of bin cluster C_i attached to two representatives of C_i , namely, bins v_k and v_l . Other bins of C_i as well as other bank bins attached to it are not displayed. $D_{kl} \leq \max_{v_p, v_q \in C} D_{pq}$.

From the inequality for $\gamma_j^{(i)}$ in the theorem's statement, $\gamma_j^{(i)} + \gamma_l^{(i)} \geq 2 \times \frac{1}{2} \max_{v_p, v_q \in C_i} D_{pq} \geq D_{kl}$, while $\gamma_j^{(i)} + D_{kl} \geq \gamma_l^{(i)}$ trivially holds. Thus, the triangle inequality holds for each triangle introduced into the ground distance by the bank bin. When extending histogram \tilde{P} to P' , the same reasoning applies to the single added bank bin and ground distance \tilde{D} , and, as a result, the triangle inequality holds for D' as well. Thus, by Theorem 1, EMD' , and, hence, EMD^* is metric. \square

5. EFFICIENT COMPUTATION OF SND

SND is defined (3) as a linear combination of several instances of EMD^* , and, thus, computation of SND involves:

- Computing the ground distance D based on the structure of the underlying network $G = \langle V, E \rangle$ ($|V| = n$, $|E| = m$) and the opinions of the users in one of the network states G_1, G_2 under comparison.
- Computing EMD^* , given that both the histograms corresponding to the network states and the ground distance have been computed.

Computation of the ground distance D implies computing the shortest paths in the network. A direct computation of D for all pairs of users using Johnson's algorithm [15] for sparse G would incur time cost $\mathcal{O}(n^2 \log n)$. Computing EMD^* is algorithmically equivalent to computing EMD , and, since the latter is formulated as a solution of a transportation problem, it can be computed either using a general-purpose linear solver, such as Karmakar's algorithm, or a solver that exploits the special structure of the transportation problem, such as the transportation simplex algorithm. The time complexity of both algorithms, however, is super-cubic in n . Thus, a precise computation of SND using existing techniques is prohibitively expensive at the scale of real-world online social networks. Furthermore, the existing approximations of EMD are either not applicable to the

comparison of histograms derived from the states of a social network, since they drastically simplify the ground distance [26, 17], or are effective only for some graphs, such as trees, structurally not characteristic of social networks [20].

We propose a method to compute SND precisely in time linear in n under the following two realistic *assumptions*.

Assumption 1: The number n_Δ of users who change their opinion between two network states G_1 and G_2 under comparison is significantly smaller than the total number n of users in the network. This assumption is reasonable, because in most applications the network states under comparison are not very far apart in time and, hence, $n_\Delta \ll n$.

Assumption 2: The opinion transportation costs, defined as the elements of adjacency matrix A^{ext} in (2), are positive integers bounded from above by constant $U \ll +\infty \in \mathbb{Z}^+$. This assumption is easy to satisfy by the appropriate choice of costs, and does not limit our analysis.

Since, according to the definition (3) of SND, its computation is equivalent to four computations of EMD^* , we will, first, focus on fast computation of EMD^* on the inputs supplied by SND. Our method for efficient computation of SND requires the following two lemmas.

LEMMA 1. *For any two histograms $P \in \mathbb{R}^n$ and $Q \in \mathbb{R}^m$ and ground distance $D \in \mathbb{R}^{n \times m}$, removal of empty bins from P and Q as well as the corresponding rows and columns from D does not affect the value of $EMD^*(P, Q, D)$.*

The proof of Lemma 1 is straightforward, since empty bins do not supply or demand any mass in the underlying transportation problem, and, hence, do not affect the cost of the optimal transportation plan. While Lemma 1 allows to remove redundant suppliers and consumers from the underlying transportation problem, the following Lemma 2 allows to transform the histograms, without affecting the value of EMD^* , exposing the redundant suppliers and consumers for removal.

LEMMA 2. *Given two arbitrary histograms $P, Q \in \mathbb{R}^n$ and a ground distance $D \in \mathbb{R}^{n \times n}$, if D is semimetric², then for any $i \in [1; n]$, the following holds*

$$EMD^*(P, Q, D) = EMD^*(P_1, \dots, P_{i-1}, P_i - \min\{P_i, Q_i\}, P_{i+1}, \dots, P_n, \\ Q_1, \dots, Q_{i-1}, Q_i - \min\{P_i, Q_i\}, Q_{i+1}, \dots, Q_n, D).$$

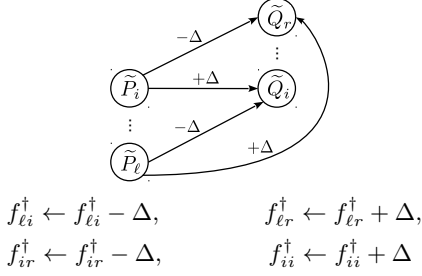
PROOF. First, we will show that there is always an optimal plan f_{ij} in the problem of optimal transportation of mass from \tilde{P} to \tilde{Q} over \tilde{D} such that $\forall i \in [1; n] : f_{ii} = \min\{\tilde{P}_i, \tilde{Q}_i\} = M$, and, then, use such a plan to argue about the value of EMD^* .

Consider an arbitrary optimal transportation plan \hat{f}_{ij} , and assume that $\exists i \in [1; n] : \delta = M - \hat{f}_{ii} > 0$. We will now use \hat{f} to construct another optimal transportation plan f_{ij}^\dagger such that $f_{ii}^\dagger = M$. Initially, we put $f^\dagger = \hat{f}$ and, then, re-route mass flows in f^\dagger to eventually achieve the desired value of f_{ii}^\dagger .

Since, initially, $f_{ii}^\dagger < M$, the remaining at least δ units of mass should be distributed by \tilde{P}_i and consumed by \tilde{Q}_i to/from other consumers/suppliers. Among those, let us pick the ones that supply/consume the least amount of mass

²Under a *semimetric* we understand a metric with symmetry requirement dropped.

to \tilde{Q}_i and from \tilde{P}_i , respectively: $\ell = \arg \min_{j \neq i} f_{ji}^\dagger$, and $r = \arg \min_{j \neq i} f_{ij}^\dagger$. W.l.o.g., let us assume that $f_{\ell i}^\dagger \leq f_{ir}^\dagger$ and denote $\Delta = \min \{f_{\ell i}^\dagger, \delta\}$. Now, we will re-route Δ units of mass in f^\dagger as follows:



The updated transportation plan is legal, as the total amount of mass supplied or consumed by each bin has not changed. The total cost of f^\dagger has been updated as follows

$$\begin{aligned} \text{newcost}(f^\dagger) &\leftarrow \text{cost}(f^\dagger) - \Delta \tilde{D}_{\ell i} - \Delta \tilde{D}_{ir} + \Delta \tilde{D}_{ii} + \Delta \tilde{D}_{\ell r} = \\ &= (\text{since } D \text{ and, hence, } \tilde{D} \text{ is semimetric, } \tilde{D}_{ii} = 0) = \\ &= \text{cost}(f^\dagger) - \Delta(\tilde{D}_{\ell i} + \tilde{D}_{ir} - \tilde{D}_{\ell r}) \leq \\ &\leq (\text{since } \tilde{D} \text{ is semimetric, } \tilde{D}_{\ell i} + \tilde{D}_{ir} \geq \tilde{D}_{\ell r}) \leq \text{cost}(f^\dagger). \end{aligned}$$

Since the cost of the obtained legal plan f^\dagger cannot be less than the cost of an optimal plan \hat{f} , the performed update of f^\dagger has not changed its cost, and the updated f^\dagger is still an optimal plan. The described above re-routing procedure is repeatedly performed on f^\dagger until f_{ii}^\dagger reaches $M = \min \{\tilde{P}_i, \tilde{Q}_i\}$.

Finally, to see why the statement of the lemma holds, we observe that the value of EMD^* is the cost of any optimal transportation plan, and the cost of f^\dagger in particular. However, the cost of f^\dagger does not depend on f_{ii}^\dagger , since, due to semimetricity of \tilde{D} , mass f_{ii}^\dagger gets transported at cost $\tilde{D}_{ii} = 0$. Thus, M can be subtracted from \tilde{P}_i , \tilde{Q}_i , and f_{ii}^\dagger , without affecting the total cost of f^\dagger . The solution of the latter modified transportation problem, however, is exactly

$$\begin{aligned} \text{EMD}^*([P_1, \dots, P_{i-1}, P_i - M, P_{i+1}, \dots, P_n], \\ [Q_1, \dots, Q_{i-1}, Q_i - M, Q_{i+1}, \dots, Q_n], D). \end{aligned}$$

□

Now, we will state our main result about the efficient computation of SND as Theorem 4, whose constructive proof will describe the computation method.

THEOREM 4. *Under Assumptions 1 and 2, SND between network states $P = [P_1, \dots, P_n]$ and $Q = [Q_1, \dots, Q_n]$ defined over network $G = \langle V, E \rangle$, ($|V| = n, |E| = m$) can be computed precisely in time*

$$\mathcal{O}(n_\Delta(m + n\sqrt{\log U} + n_\Delta^2 \log(n_\Delta nU))).$$

PROOF. We will focus on the efficient computation of the first summand $\text{EMD}^*(P^+, Q^+, D(P, +))$ in the definition (3) of $\text{SND}(P, Q, D)$, as computation of three other summands is procedurally equivalent and takes the same time. For the analysis of the computation of $\text{EMD}^*(P^+, Q^+, D(P, +))$, let us assume, without loss of generality, that $\sum_{i=1}^n P_i^+ \geq \sum_{j=1}^n Q_j^+$. Let us also assume, for the ease of explanation, that the histograms are extended with one bank per bin. By definition (4), $\text{EMD}^*(P^+, Q^+, D(P, +))$ is the solution of a

transportation problem with supplies $\tilde{P}^+ = [P_1^+, \dots, P_n^+, 0, \dots, 0]$, demands $\tilde{Q}^+ = [Q_1^+, \dots, Q_n^+, B_1, \dots, B_n]$ and ground distance $\tilde{D}(P, +)$, where B_i is the bank bin attached to bin Q_i , and the histograms and the ground distance extended according to the definition (4) of EMD^* .

Now, we can apply Lemmas 1 and 2 to reduce the size of the obtained transportation problem. From Assumption 2, $\tilde{D}(P, +)$ is semimetric. Non-negativity and identity of indiscernibles straightforwardly follow from Assumption 2 and the definition of the length of a shortest path. Subadditivity follows from the shortest path problem's optimal substructure. Thus, we can apply Lemma 2 to each pair $\tilde{P}_i^+, \tilde{Q}_i^+$ of corresponding suppliers and consumers, and due to Assumption 1, a large number $(n - n_\Delta)$ of them have equal values. As a result, many suppliers and consumers become empty. Then, due to Lemma 1, all the obtained empty bins can be disregarded. If we put $M_i = \min \{P_i^+, Q_i^+\}$, then the reduced transportation problem is defined for suppliers $[P_{i_1}^+ - M_{i_1}, \dots, P_{i_{n_\Delta}}^+ - M_{i_{n_\Delta}}]$ and consumers $[Q_{j_1}^+ - M_{j_1}, \dots, Q_{j_{n_\Delta}}^+ - M_{j_{n_\Delta}}, B_1, \dots, B_n]$, and ground distance $\tilde{D}(P, +)$ that contains only the rows and columns corresponding to the remaining suppliers and consumers. The remaining suppliers and non-bank consumers are those that correspond to the users who have different opinion in P^+ and Q^+ , and the number of such users, due to Assumption 1, is at most n_Δ . The bank bins, however, do not get affected by Lemma 2 (since only the banks of the lighter histogram can have non-zero mass) in \tilde{Q}^+ and hence do not get removed, yet they get removed from \tilde{P}^+ due to Lemma 1 as being empty. Thus, we have ended up with an unbalanced transportation problem, where the number n_Δ of suppliers is much less than the number $n + n_\Delta$ of consumers.

Now, in order to compute $\text{EMD}^*(P^+, Q^+, D(P, +))$, we need to compute $\tilde{D}(P, +)$ and to actually solve the obtained transportation problem.

Due to the structure of the reduced transportation problem, we need to compute only a small part of $\tilde{D}(P, +)$. Since there are at most n_Δ suppliers, we need to solve at most n_Δ instances of single-source shortest path problem with at most $n_\Delta + n$ destinations. Since, due to Assumption 2, edge costs in the network are integer and bounded by U , each instance of a single-source shortest path problem can be solved using Dijkstra's algorithm based on a combination of a radix and a Fibonacci heaps [1] in time

$$T_{sssp} = \mathcal{O}(m + n \log \sqrt{U}).$$

(Notice, that if we assumed $\sum_{i=1}^n P_i^+ \leq \sum_{j=1}^n Q_j^+$, and the reduced \tilde{P}^+ contained $n_\Delta + n$ bins, we would not need to run $n_\Delta + n$ instances of Dijkstra's algorithm. Instead, we would invert the edges in the network and compute shortest paths in reverse, still performing only n_Δ single-source shortest path computations.)

Next we approach the solution of the reduced transportation problem with known ground distances. This problem can be viewed as a min-cost network flow problem in an unbalanced bipartite graph. Since, due to Assumption 2, edge costs are integers bounded by U , our min-cost flow problem can be solved using Goldberg-Tarjan's algorithm [11]

augmented with the two-edge push rule of [2] in time

$$T_{transp} = \mathcal{O}(n_{\Delta}m + n_{\Delta}^3 \log(n_{\Delta} \max_{i,j} D(P, +)_{ij})).$$

Since no shortest path has more than $(n - 1)$ edge, and the edge costs are bounded by U , the expression for time simplifies to

$$T_{transp} = \mathcal{O}(n_{\Delta}m + n_{\Delta}^3 \log(n_{\Delta}nU)).$$

Thus, the total time for computing $\text{EMD}(P^+, Q^+, D(P, +))$ and, consequently, $\text{SND}(P, Q, D)$ is

$$T = \mathcal{O}(n_{\Delta}T_{ssp} + T_{transp}) = \mathcal{O}(n_{\Delta}(m + n \log \sqrt{U} + n_{\Delta}^2 \log(n_{\Delta}nU))).$$

□

Notice that, if the social network is sparse, that is $m = \mathcal{O}(n)$, and the number of changes n_{Δ} is bounded, then, according to Theorem 4, SND is computable in time $\mathcal{O}(n)$.

6. EXPERIMENTAL RESULTS

In this section, we report experimental results, demonstrating the utility of SND in applications and comparing it with other distance measures. We also study the scalability of our implementation of SND.

6.1 Experimental Setup

Real-World Data: Our Twitter dataset, based on data from [19] contains 10k users, each having an average of 130 follower-followee edges. Among the tweets sent by these users between May-2008 and August-2011, we select those relevant to political topics, such as “Obama”, “GOP”, “Palin”, “Romney”. We break the entire observation period into quarters and, in each quarter, assess every user’s polarity with respect to our topics of interest. Polarity of all users within one quarter comprise a network state.

Synthetic Data: We also perform experiments on synthetic scale-free networks of sizes $|V|$ from 10k to 200k and scale-free exponents from -2.9 to -2.1 . To generate the first network state, a number of initial adopters are chosen uniformly at random, and approximately equal numbers of them adopt “+” and “−” opinions. Each subsequent network state G_{i+1} is randomly generated from the preceding network state G_i as follows. A number of G_i ’s neutral users get a chance to be activated. Each of them adopts an opinion from her neighbors with probability \mathbb{P}_{nbr} and a random opinion with probability \mathbb{P}_{ext} . If a user is to adopt an opinion from the neighbors, which opinion to adopt is decided in a probabilistic voting fashion based on the numbers of active in-neighbors of each kind.

Distance Measures: We compare SND with the following distance measures.

- *hamming*(P, Q). The Hamming distance is a representative of all the distance measures performing coordinate-wise comparison.

- *quad-form*(P, Q, L) = $\sqrt{(P - Q)L(P - Q)^T}$. Quadratic-Form Distance [14] based on the Laplacian L [22] of the network. It takes the differences of opinions of the corresponding users and combines them based on the network’s structure

- *walk-dist*(P, Q) = $\frac{1}{n} \|cnt(P) - cnt(Q)\|_1$. Compares vectors $cnt(P) = [cnt(P_1), \dots, cnt(P_n)]$ of users’ “contention”, where $cnt(P_i)$ is the amount by which the i ’th user’s opinion deviates from the opinion of this user’s average active in-neighbor. Thus, *walk-dist* summarizes how different the network’s users are from their respective neighbors.

6.2 Detecting Anomalous Network States

Synthetic Data: In a series of network states, we want to detect which ones are anomalous. In particular, we are interested in the anomalies which are hard to detect by observing the summary of the social network (e.g., the number of new activations). Thus, in experiments with synthetic data, to simulate an anomaly, we change the values of \mathbb{P}_{nbr} and \mathbb{P}_{ext} preserving their sum, thereby, affecting only qualitatively the process of new users’ activation. In a series of network states, we compute the distances between the adjacent states, normalize these distances by the number of active users, and scale. Then, spikes in the resulting series of distances are considered anomalies.

A qualitative analysis of anomaly detection on synthetic data is presented in Fig. 7. For each simulated anomaly,

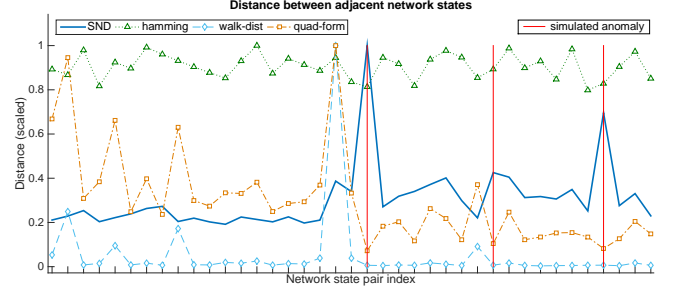


Figure 7: Anomaly detection on synthetic data. $|V| = 20k$, scale-free exponent $\gamma = -2.3$. A series of 40 network states is generated using $\mathbb{P}_{nbr} = 0.12$ and $\mathbb{P}_{ext} = 0.01$ for normal and $\mathbb{P}_{nbr} = 0.08$ and $\mathbb{P}_{ext} = 0.05$ for anomalous network states’ generation.

SND produces a well noticeable spike, while other distance measures do not recognize such anomalies.

In order to quantify the performance of the competing distance measures at detecting the true simulated anomalies, we create a simple anomaly score $S_t = (d_t - d_{t-1}) + (d_t - d_{t+1})$, where d_t is the value of a given distance measure at time t normalized by the number of users active at time t and scaled. We rank the network state transitions for each compared distance measure by S_t in decreasing order and compute true and false positive rates for increasing ranks. The corresponding ROC curves are displayed in Fig. 8.

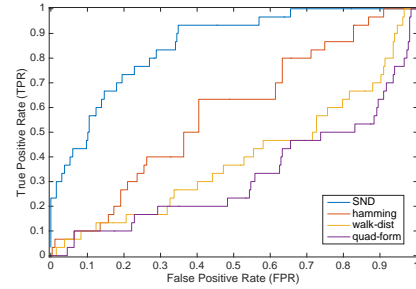


Figure 8: ROC curves comparing the quality of anomaly detection by different distance measures in a series of 300 network states over synthetic network with $|V| = 30k$ and scale-free exponent $\gamma = -2.3$. The network states are generated using $\mathbb{P}_{nbr} = 0.08$ and $\mathbb{P}_{ext} = 0.001$ for normal and $\mathbb{P}_{ext} = 0.011$ and $\mathbb{P}_{nbr} = 0.07$ for anomalous instances.

SND’s accuracy dominates that of competing distance measures throughout the spectrum of false positive rates. Particularly, for false positive rates up to 0.3, SND achieves a

true positive rate of 0.83, while the next best distance measure (*hamming*) achieves only 0.4.

Twitter Data: To obtain the ground truth for anomaly detection on our Twitter dataset, we collect “search interest” data from Google Trends³ and cross-check this data with American Presidents⁴ log of important political events in the US. The anomaly detection results for topic “Obama” are shown in Fig. 9.

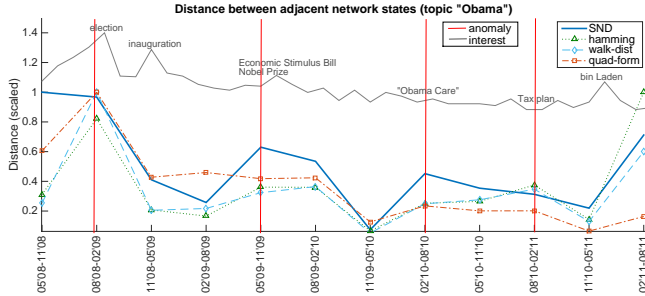


Figure 9: Anomaly detection on Twitter data (May’08-Aug’11) for topic “Obama”. The plots for distances between network states are accompanied by the plot showing Google Trends’ (scaled) interest in topic “Obama”. Network states detected to be anomalous by at least one distance measure are indicated with red vertical lines.

We can distinguish two types of network states and, hence, events based on SND’s behavior relatively to other distance measures. One type is the events corresponding to network states where SND agrees with other distance measures. Two examples are (i) the first anomaly – Barack Obama’s election for the President of the US, and (ii) bin Laden’s death being the last spike on the Google Trends curve (even though, all distance measures noticeably increase their value during the last quarter, we do not mark this quarter as anomalous, since we do not have the distance values for the next quarter.) These events are unlikely to have been perceived differently by the US users of Twitter and, hence, probably have not provoked a polarized response.

The other type of events are those where SND noticeably disagrees with other distance measures. For example, during quarters 05’09-11’09, the Economic Stimulus Bill had a highly polarized response in the House of Representatives⁵, with no Republican voting in its favor. Another such anomaly takes place during quarters 02’10-08’10, when the Affordable Care Act (“Obama Care”) was introduced, and which was and still remains a very controversial topic. The latter can be seen from the House vote distribution⁶ and from the fact that, according to socialmention.com⁷, even in October 2015, “Obama Care” is still perceived equally positively and negatively in microblogs.

6.3 Predicting User Opinions

Given a series of states of a social network we want to predict the unknown opinions of the users in the current network state G_0 based on the observed recent G_{-t} ($t \in \mathbb{N}$) and the (incomplete) current network states. For example,

³<http://www.google.com/trends/explore>

⁴<http://www.american-presidents-history.com>

⁵<http://www.nytimes.com/2009/01/29/us/politics/29obama.html>

⁶Democrats – 219 yeas, republicans – 212 nays (<http://www.healthreformvotes.org/congress/roll-call-votes/h165-111.2010>)

⁷<http://socialmention.com/search?q=%22obama+care%22>

if certain Twitter users have not tweeted (enough) in the current quarter, we may want to predict the opinions of these users in the current quarter based on the observed opinions of all users of the network. We assume that during the periods corresponding to the observed recent network states G_{-t} , the social network evolved “smoothly”, that is, the recent states can predict the current state. Under this assumption, we use a distance measure to compute the distances $\text{dist}(G_{-t}, G_{-t+1})$ between the adjacent past network states, then, extrapolate the obtained series to estimate the distance d^* from the most recent G_{-1} to the yet unknown *complete* current network state. Then, we assign different opinions to the target users in the current network state, trying to make the distance $\text{dist}(G_{-1}, G_0^*)$ from the most recent to the modified current network state as close to estimate d^* as possible. The search for the best assignment of opinions to the target users is randomized — the number of the uniformly randomly generated opinion assignments for all target users is considerably lower than the total number of possible assignments (we use 100 random opinion assignments in each experiment). In each experiment, we uniformly randomly select 20 active users – with approximately equal number of positive and negative users – in the current network state, predict their opinions and measure the prediction accuracy. This procedure is repeated 10 times, and mean accuracies and standard deviations are reported.

The predictions are made using the above distance-based method with SND as well as other distance measures. To put the prediction performance of these methods in context, we add into comparison two non-distance-based methods – one basic, and one state-of-the-art – that make predictions based on the known quantified opinions of the users and the network’s structure. One such method, *nhood-voting*, derives the opinion of each target user based on the opinions of this user’s active in-neighbors in a probabilistic voting fashion, or selects it uniformly randomly in the absence of active in-neighbors. Another method, *community-lp* [9, IV.B] detects communities in the network via label propagation and, then, predicts user opinions based on these users’ membership in the discovered communities.

We experiment on both synthetic and real-world data. For synthetic data, we generate a scale-free network with $n = 10k$ users and scale-free exponent $\gamma = -2.5$. A series of network states is generated using the same algorithm as in the case with anomaly detection, with probabilities of opinion adoption from the neighborhood and from the “external source” ranging between 0.001 and 0.2. The number of initial adopters in the first network state is 800. We use 3 most recent network states to estimate the distance from the most recent to the incomplete current network state.

Results for opinion prediction are summarized in Table 1. There are three important observations:

- Firstly, among the distance-based methods, the one that uses SND always performs best, with an average prediction accuracy of 74-75% and a consistently low standard deviation. This suggests that SND captures more opinion dynamics-specific information than the other distance measures, and should be preferred, particularly, when such simple statistics as the rate of new user activation are uninformative.

- Secondly, SND-based prediction method works considerably better than method *nhood-voting* that bases the opinion prediction for each user on the opinions of the user’s

in-neighbors. This emphasizes the importance of analyzing opinion dynamics at the level of the entire network, rather than for each egonet in isolation.

- Lastly, SND-based method outperforms the state-of-the-art method *community-lp*, based on community detection via label propagation. In our experiments, *community-lp*'s prediction accuracy is 57-65%, while this method's authors report the accuracy of 95% for their data [9]. The likely cause of such a discrepancy is a *very* high level of homophily in their data (users almost exclusively follow those having the same opinion), while in our less homophilous data, *community-lp* performs worse by capturing only users' reachability by the opinions of each kind, and SND performs better by looking for the *most probable* opinion propagation scenario.

| User Opinion Prediction Accuracy, % | | | | |
|-------------------------------------|----------------|-------------|-----------------|-------------|
| Method | Synthetic Data | | Real-World Data | |
| | μ | σ | μ | σ |
| SND | 74.33 | 2.65 | 75.63 | 5.60 |
| hamming | 68.44 | 12.34 | 68.13 | 5.80 |
| quad-form | 66.67 | 13.58 | 67.50 | 9.63 |
| walk-dist | 56.22 | 15.35 | 31.88 | 9.98 |
| nhood-voting | 62.11 | 8.58 | 61.25 | 5.82 |
| community-lp | 65.25 | 9.43 | 56.87 | 8.43 |

Table 1: Means μ and standard deviations σ of user opinion prediction accuracies.

6.4 Sensitivity to Opinion Dynamics Models

In this section, we show the effectiveness of SND in detecting qualitative changes in the network's evolution w.r.t. an advanced opinion dynamics model, that cannot be spotted by the distance measures performing coordinate-wise comparison. We generate a number of pairs $\langle G_1, G_2 \rangle$ of subsequent network states over a synthetic scale-free network. Some of these pairs correspond to *normal transitions*, while others correspond to *anomalous transitions* in the network's evolution. For the normal transitions, G_2 is generated from G_1 using the Independent Cascade Model with Competition (ICC) [7]. For anomalous transitions, most new activations in G_2 happen randomly, not relying on the network's structure. We study the distances assigned to normal and anomalous transitions by SND and ℓ_1 , and plot them as functions of the number n_Δ of users having different opinion in G_1 and G_2 of each transition. According to the results in Fig. 10, SND clearly separates anomalous transitions from normal ones. ℓ_1 , however, cannot discern anomalous transitions, as its value is mostly determined by n_Δ , which is representative of the distance measures performing coordinate-wise comparison.

6.5 Scalability of SND

We implemented⁸ SND in MATLAB, with parts written in C++. We use a min-cost network flow solver CS2 [12] that implements Goldberg-Tarjan's algorithm [11], but, unlike it is prescribed by Theorem 4, does not use the two-edge push rule of [2]. Additionally, for computing shortest paths, our implementation of Dijkstra's algorithm uses a priority queue based on a binary heap, rather than a combination

⁸<https://cs.ucsb.edu/~victor/pub/ucsb/dbl/snd/>

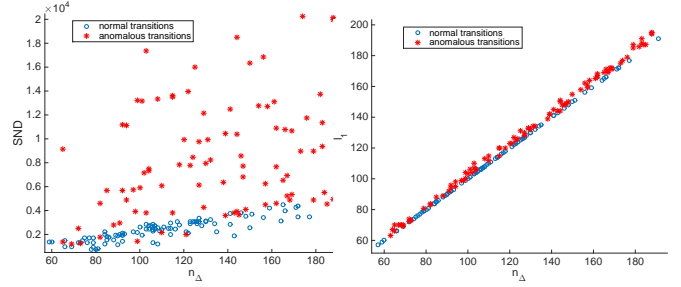


Figure 10: SND and ℓ_1 distances between network states of normal (ICC) and anomalous (random) transitions.

of a Fibonacci and a radix heaps [1]. As a result, our implementation of SND scales slightly worse than linearly as guaranteed by Theorem 4, but still very well to be applicable to real-world social networks. Fig. 11 shows how our implementation of SND based on the proposed efficient method scales in the number n of users in the network in comparison with a direct computation of SND based on a CPLEX linear solver. Our implementation's scalability in the number n_Δ of users who have changed their opinion is shown in Fig. 12.

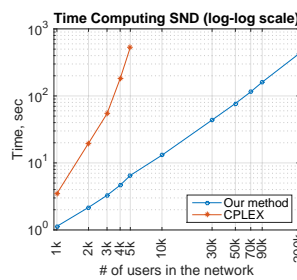


Figure 11: Time for computing SND when the number of users having different opinion is fixed at $n_\Delta = 1000$ and the total number of users n in the network grows up to 200k.

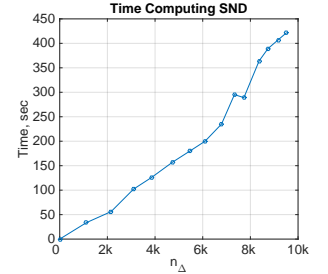


Figure 12: Time for computing SND using our method when the network size is fixed at $n = 20k$, and the number of users n_Δ having different opinions grows up to 10k.

7. RELATED WORK

There is a large number of existing distance measures used in vector spaces, including ℓ_p , Hamming, Canberra, Cosine, Kullback-Leibler, and Quadratic Form [14] distances. However, none of them is adequate for the comparison of network states, since these distance measures either compare vectors coordinate-wise, thereby, not capturing the interaction between users in the network, or in the case of Quadratic Form distance, capture the user interaction in a very limited and uninterpretable way.

Existing graph-oriented distance measures are also unsuitable for comparing network states with polar opinions. The first class of such distance measures is graph isomorphism-based distance measures, such as *largest common subgraph* [6]. These distance measures are node state-oblivious, and, hence are not applicable to the comparison of network states. Another class of graph distance measures is *Graph Edit Distance (GED)-based* measures [10] that define the distance

between two networks as the cost of the optimal sequence of edit operations, such as node or edge insertion, deletion, or substitution, transforming one network into another. GED can be node state-aware, but its value is not interpretable from the opinion dynamics point of view, and even its approximate computation takes time cubic in $|V|$ (a single computation of GED on a 10k-node network on our hardware takes about a month).

A third class of distance measures includes *iterative distance measures* [4, 16, 21], which express similarity of the nodes of two networks recursively, use a fix-point iteration to compute node similarities, and, then, aggregate node similarities to obtain the similarity of two networks. Iterative distance measures share the problem of GED – they do not capture the way opinions spread in the network.

The last class includes *feature-based* distance measures [3, 27, 29], which compare either the distributions of local node properties (e.g., degree, clustering coefficient) or the spectra of two networks. Despite their efficient computability, such distance measures do not fit the comparison of network states with polar opinions. The spectral distance measures are inadequate because they do not deal with node states directly⁹, while other feature-based distance measures only deal with summaries based on opinion of each kind and, thus, cannot capture the competition of polar opinions.

8. LIMITATIONS

Despite the demonstrated effectiveness and efficiency of SND, there are scenarios in which its use is either prohibitively or unnecessarily expensive.

- One reason to choose a simpler distance measure, such as ℓ_p , over SND is the latter’s computational cost. While it is asymptotically linear in the number of nodes, it can, potentially, be too high in practice for networks having 100M+ nodes, where a single computation of SND can take several days. If the use of a simpler distance measure is undesirable, one can partition the network into clusters of tractable size and perform the SND-based analysis on each cluster.

- Another scenario when using SND may be excessive is when the changes in the rate of new user activation reveal enough information for the target application (for example, the activation rate alone is clearly enough to detect the US presidential election day), and, in such a case, the distance measures as simple as Hamming distance may suffice.

9. FUTURE RESEARCH

Among the directions for future research are the following.

- Since SND is, effectively, the first distance measure designed specifically for the comparison of states of a social network containing competing opinions, one potential future research direction is using SND in other applications operating in a metric space setting, such as network state classification, clustering, and search.

- Additionally, it may be lucrative to combine SND with non-distance-based methods. Thus, in the method of [9] that predicts opinions based on the content of the users’ tweets, the objective function can be augmented with an SND-based term, thereby, performing opinion fitting at both the micro-level of each user and the macro-level of the entire network.

⁹Even if node states are artificially encoded into a network’s structure, there is still a possibility for two structurally different networks to have identical spectra and, hence, a zero spectral distance.

- Finally, it may be fruitful to design a distance measure that would capture changes in both the opinions of the users and the structure of the social network simultaneously. Such a distance measure would be more computationally complex than SND due to the network alignment requirement, yet, useful for the comparison of network states defined over very different networks.

10. CONCLUSION

In this paper, we proposed Social Network Distance (SND) – the first distance measure for comparing the states of a social network containing competing opinions. Our distance measure quantifies how likely it is that one state of a social network has evolved into another state under a given model of polar opinion propagation. Despite the high computational complexity of the transportation problem underlying SND, we propose a linear-time algorithm for its precise computation, making SND applicable to real-world online social networks. We demonstrate the usefulness of SND in detecting anomalous network states and predicting user opinions in both synthetic and real-world data, where it consistently outperforms other distance measures. Our anomaly detection method achieves a true positive rate (TPR) of 0.83, while the next best method’s TPR is only 0.4. The accuracy of SND-based method for user opinion prediction averages at 75.63%, which is 7.5% higher than that of the next best method. We also show that, unlike the distance measures performing coordinate-wise comparison, SND can detect qualitative changes in the network’s evolution pattern.

Our results emphasize the importance of taking into account user locations in the analysis of social networks, and that the analysis of opinion dynamics at the level of an entire social network can provide more information about the network’s evolution than the methods operating at the level of egonets.

11. REFERENCES

- [1] R. K. Ahuja, K. Mehlhorn, J. Orlin, and R. E. Tarjan. Faster algorithms for the shortest path problem. *Journal of the ACM (JACM)*, 37(2):213–223, 1990.
- [2] R. K. Ahuja, J. B. Orlin, C. Stein, and R. E. Tarjan. Improved algorithms for bipartite network flow. *SIAM Journal on Computing*, 23(5):906–933, 1994.
- [3] M. Berlingerio, D. Koutra, T. Eliassi-Rad, and C. Faloutsos. NetSimile: a scalable approach to size-independent network similarity. *arXiv preprint arXiv:1209.2684*, 2012.
- [4] V. D. Blondel, A. Gajardo, M. Heymans, P. Senellart, and P. Van Dooren. A measure of similarity between graph vertices: Applications to synonym extraction and web searching. *SIAM review*, 46(4):647–666, 2004.
- [5] A. Borodin, Y. Filmus, and J. Oren. Threshold models for competitive influence in social networks. In *Internet and Network Economics*, pages 539–550. Springer, 2010.
- [6] H. Bunke and K. Shearer. A graph distance metric based on the maximal common subgraph. *Pattern recognition letters*, 19(3):255–259, 1998.
- [7] T. Carnes, C. Nagarajan, S. M. Wild, and A. Van Zuylen. Maximizing influence in a competitive

- social network: a follower's perspective. *EC*, pages 351–360, 2007.
- [8] K. L. Clarkson. Nearest-neighbor searching and metric space dimensions. *Nearest-neighbor methods for learning and vision: theory and practice*, pages 15–59, 2006.
 - [9] M. D. Conover, B. Gonçalves, J. Ratkiewicz, A. Flammini, and F. Menczer. Predicting the political alignment of Twitter users. In *SocialCom*. IEEE, 2011.
 - [10] X. Gao, B. Xiao, D. Tao, and X. Li. A survey of Graph Edit Distance. *Pattern Analysis and applications*, 13(1):113–129, 2010.
 - [11] A. Goldberg and R. Tarjan. Solving minimum-cost flow problems by successive approximation. In *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, pages 7–18. ACM, 1987.
 - [12] A. V. Goldberg. An efficient implementation of a scaling minimum-cost flow algorithm. *Journal of algorithms*, 22(1):1–29, 1997.
 - [13] A. Goyal, F. Bonchi, and L. V. Lakshmanan. Learning influence probabilities in social networks. *WSDM*, pages 241–250, 2010.
 - [14] J. Hafner, H. S. Sawhney, W. Equitz, M. Flickner, and W. Niblack. Efficient color histogram indexing for quadratic form distance functions. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 17(7):729–736, 1995.
 - [15] D. B. Johnson. Efficient algorithms for shortest paths in sparse networks. *Journal of the ACM (JACM)*, 24(1):1–13, 1977.
 - [16] E. Leicht, P. Holme, and M. Newman. Vertex similarity in networks. *Physical Review E*, 73(2):026120, 2006.
 - [17] L. Li, M. Ma, P. Lei, X. Wang, and X. Chen. A linear approximate algorithm for Earth Mover's Distance with thresholded ground distance. *Mathematical Problems in Engineering*, 2014.
 - [18] V. Ljosa, A. Bhattacharya, and A. K. Singh. Indexing spatially sensitive distance measures using multi-resolution lower bounds. *EDBT*, pages 865–883, 2006.
 - [19] K. Macropol, P. Bogdanov, A. K. Singh, L. Petzold, and X. Yan. I act, therefore I judge: Network sentiment dynamics based on user activity change. *ACM ASONAM*, pages 396–402, 2013.
 - [20] A. McGregor and D. Stubbs. Sketching earth-mover distance on graph metrics. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 274–286. Springer, 2013.
 - [21] S. Melnik, H. Garcia-Molina, and E. Rahm. Similarity flooding: A versatile graph matching algorithm and its application to schema matching. In *Data Engineering, 2002. Proceedings.*, pages 117–128. IEEE, 2002.
 - [22] R. Merris. Laplacian matrices of graphs: a survey. *Linear algebra and its applications*, 197:143–176, 1994.
 - [23] S. A. Myers and J. Leskovec. The bursty dynamics of the Twitter information network. In *Proceedings of the 23rd international conference on World wide web*, pages 913–924. ACM, 2014.
 - [24] O. Pele and M. Werman. A linear time histogram metric for improved SIFT matching. In *Computer Vision-ECCV 2008*, pages 495–508. Springer, 2008.
 - [25] Y. Rubner, C. Tomasi, and L. J. Guibas. The Earth Mover's Distance as a metric for image retrieval. *International Journal of Computer Vision*, 40(2):99–121, 2000.
 - [26] Y. Tang, U. Leong Hou, Y. Cai, N. Mamoulis, and R. Cheng. Earth Mover's Distance based similarity search at scale. *Proceedings of the VLDB Endowment*, 7(4):313–324, 2013.
 - [27] R. C. Wilson, E. R. Hancock, and B. Luo. Pattern vectors from algebraic graph theory. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 27(7):1112–1124, 2005.
 - [28] E. Yildiz, D. Acemoglu, A. E. Ozdaglar, A. Saberi, and A. Scaglione. Discrete opinion dynamics with stubborn agents. *Available at SSRN 1744113*, 2011.
 - [29] P. Zhu and R. C. Wilson. A study of graph spectra for comparing graphs. In *BMVC*, 2005.